

Strategies for Acceleration of Deep Learning in Algorithm & CMOS-based Hardware

- V. Sze et. al., “Efficient Processing of Deep Neural Networks: A Tutorial and Survey”. Proc. IEEE Vol. 105, No. 12, Dec. 2017
- W. Dally, “High-Performance Hardware for Machine Learning”, NIPS 2015
- S. Han & W. Dally, “Deep Learning Tutorial and Recent Trends”, FPGA 2017
- A. Krizhevsky, “One weird trick for parallelizing convolutional neural networks”, arXiv, 2014
- J. Dean, “Machine Learning for Systems and Systems for Machine Learning”, NIPS 2017
- Etc..

Jan. 10. 2019
Jong Hoon Shin

Why Deep Learning Did Not Work in the Past?

Geoffrey Hinton summarized the findings up to today in these four points:

1. Our **labeled datasets** were thousands of times too small.

Youtube comments , reCAPTCHA, Facebook privacy, Amazon reviews..

2. Our **computers** were millions of times too slow.

GPU, FPGA, ASIC, Neuromorphic (Spiking, NVM)..

3. We **initialized the weights** in a stupid way.

RBM, Xavier/Glorot/He initialization ..

4. We used the wrong type of **non-linearity**.

ReLu series..

Advances in Deep Learning Algorithms

Non-Linearities

Relu
Sigmoid
Tanh
GRU
LSTM
Linear
...

Optimizer

SGD
Momentum
RMSProp
Adagrad
Adam
Second Order (KFac)
...

Connectivity Pattern

Fully connected
Convolutional (Conv, MaxPooling)
Dilated
Recurrent
Recursive
Skip / Residual
Random

Loss

Cross Entropy
Adversarial
Variational
Max. Likelihood
Sparse
L2 Reg
REINFORCE
...

Hyper Parameters

Learning Rate
Decay
Layer Size
Batch Size
Dropout Rate
Weight init
Data augmentation
Gradient Clipping
Beta
Momentum

Advances in Tools for Deep Learning

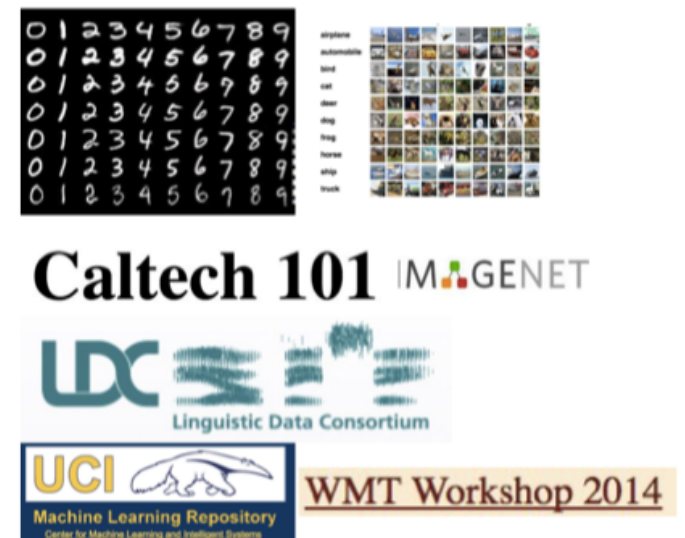
Platforms



Frameworks

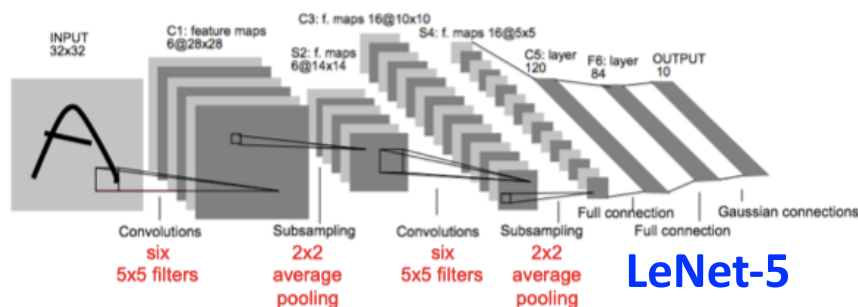


Datasets



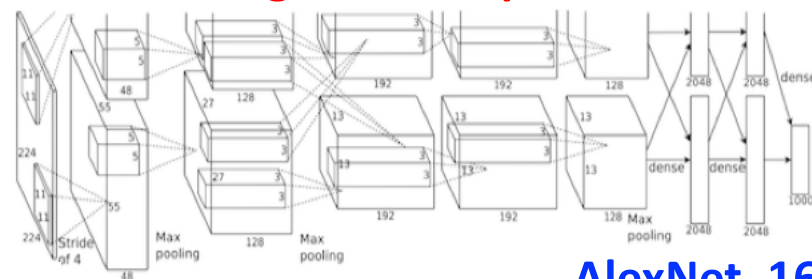
Famous DL Algorithms for MNSIT & ImageNet

MNIST

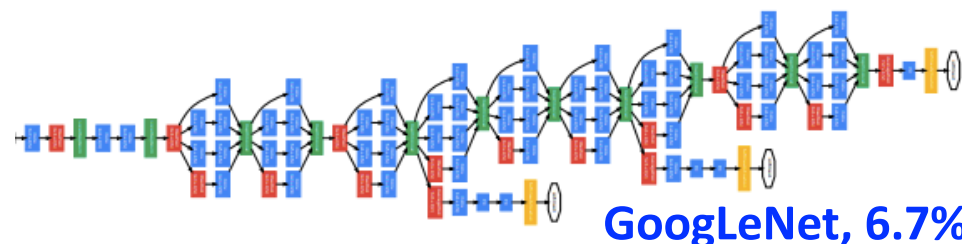


LeNet-5

ImageNet Competition

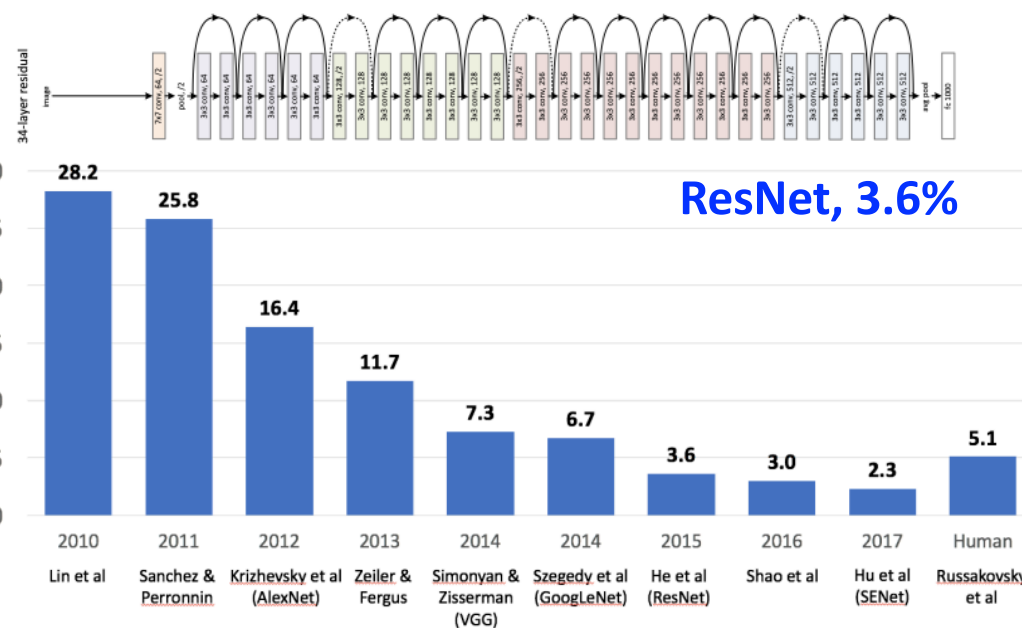


AlexNet, 16.4%



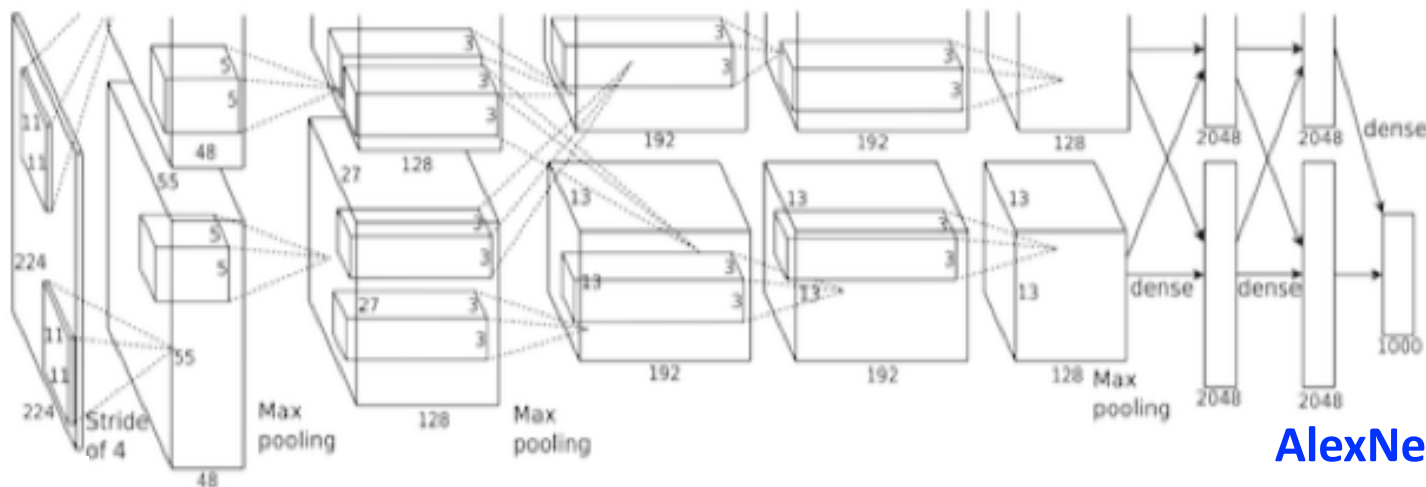
GoogLeNet, 6.7%

- MNIST → 1998'
- ImageNet → 2018'
- Sequential data & Video(30fps)?
 - Throughput matters
 - # of width & layers
 - IOT/Edge → Energy matters



ResNet, 3.6%

AlexNet's Parameters and MACs



AlexNet, 16.4%

Layer	Filter Size (RxS)	# Filters (M)	# Channels (C)	Stride
1	11x11	96	3	4
2	5x5	256	48	1
3	3x3	384	256	1
4	3x3	384	192	1
5	3x3	256	192	1

Input Image: (224x224x3)

1st Output:
 $(224-11)/4+1 \rightarrow 55 \times 55$

1st Params:
 $(11 \times 11 \times 3) \times 96 \rightarrow 34K$

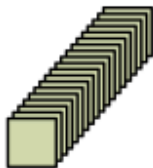
1st MACs:
 $(34K \times 55 \times 55) \rightarrow 105M$

Layer 1

34k Params
105M MACs

Layer 2

307k Params
224M MACs

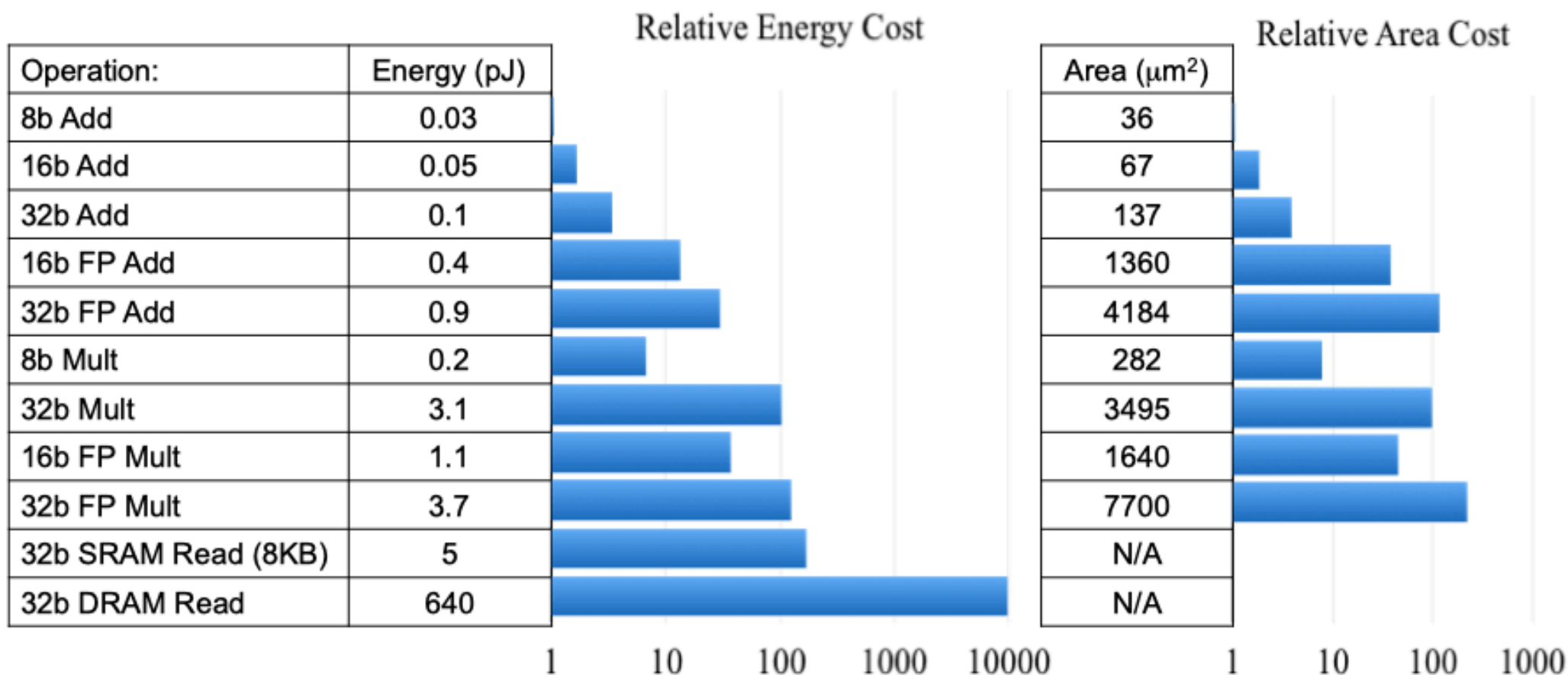
Layer 3

885k Params
150M MACs

Summary of Parameters and MACs

Metrics	LeNet-5	AlexNet	VGG-16	GoogLeNet (v1)	ResNet-50
Top-5 error	n/a	16.4	7.4	6.7	5.3
Input Size	28x28	227x227	224x224	224x224	224x224
# of CONV Layers	2	5	16	21 (depth)	49
Filter Sizes	5	3, 5, 11	3	1, 3, 5, 7	1, 3, 7
# of Channels	1, 6	3 - 256	3 - 512	3 - 1024	3 - 2048
# of Filters	6, 16	96 - 384	64 - 512	64 - 384	64 - 2048
Stride	1	1, 4	1	1, 2	1, 2
# of Weights	2.6k	2.3M	14.7M	6.0M	23.5M
# of MACs	283k	666M	15.3G	1.43G	3.86G
# of FC layers	2	3	3	1	1
# of Weights	58k	58.6M	124M	1M	2M
# of MACs	58k	58.6M	124M	1M	2M
Total Weights	60k	61M	138M	7M	25.5M
Total MACs	341k	724M	15.5G	1.43G	3.9G

- Multiply 1.2M training image for a single training epoch!
- Most of Parameters → FC & Most of MACs → Filter

Considering Energy Consumption..



- Reduce **Memory Access**
- Reduce **Data Type** : F→I, 32→16→8
- Reduce **Mult & Add**

Strategies for DL Acceleration

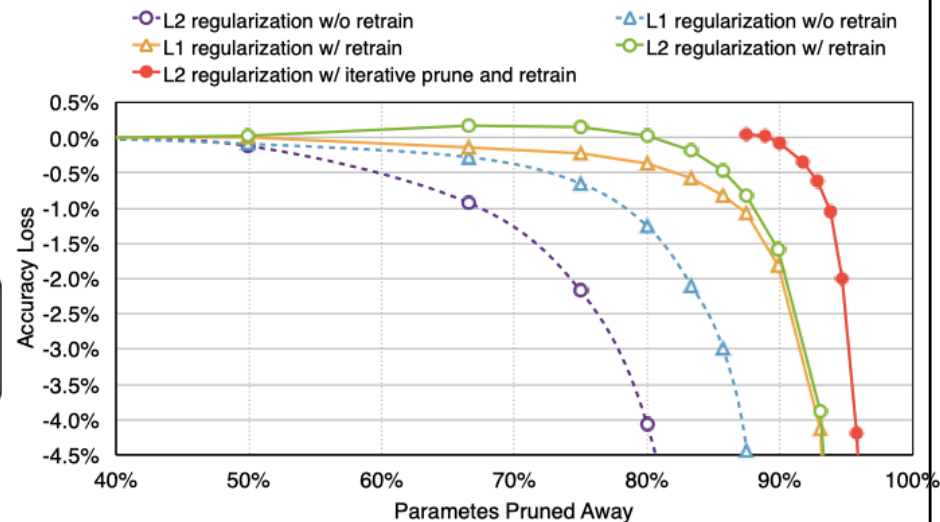
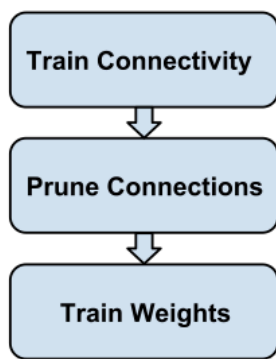
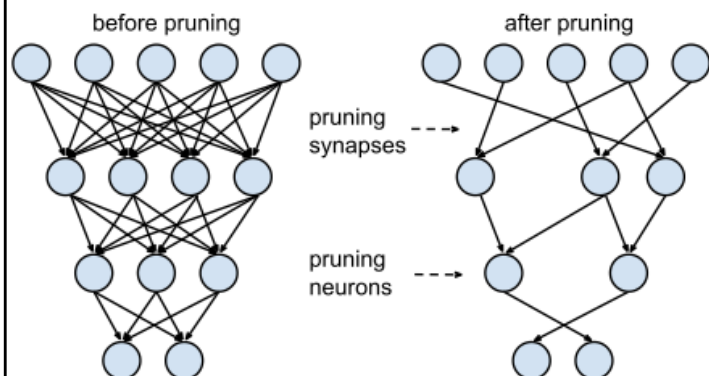
	Inference	Training
Algorithm	<p>Pruning : P-Threshold, +Training</p> <p>Weight Sharing:</p> <ul style="list-style-type: none"> w-KNN → Quant, Huffman coding <p>Quantization: Floating P. → Fixed P.</p> <p>Low Rank Approximation</p> <p>Binary/Ternary Net</p> <p>Computational Transform</p> <ul style="list-style-type: none"> Winograd, FFT, Strassen 	<p>Parallelization:</p> <ul style="list-style-type: none"> Data/Parameter/Hyp.Par <p>Mixed Precision MAC: FP16 + FP32</p> <p>Model Distillation:</p> <ul style="list-style-type: none"> Teachers(GGN,VGG,Res)&Student <p>DSD(Dense-Sparse-Dense) Training</p> <p>Mixed Precision Net:</p> <ul style="list-style-type: none"> BinaryConnect, QNN
Hardware	<p>Architecture:</p> <ul style="list-style-type: none"> Temporal(SIMD,SIMT with Global Buffer) : General purpose GPUs Spatial(Local Memory in Dist. Processing Element) <p>Data Flow Taxonomy:</p> <ul style="list-style-type: none"> Weight Stationary, Output Stationary No local reuse (M-Bottleneck: Roofline) Row stationary <p>Data type: Int8/16 for inference(TPU1), Mix FP16/32 for training(TPU2)</p> <p>Data Reuse: Conv, Feature map, Filter</p> <p>Kernel Computation(Toeplitz Matrix, TensorCore, Matrix-Multiply Unit)</p> <p>Sparse, Compressed Model(EIE, NVIDIA)</p>	

Strategies for DL Acceleration

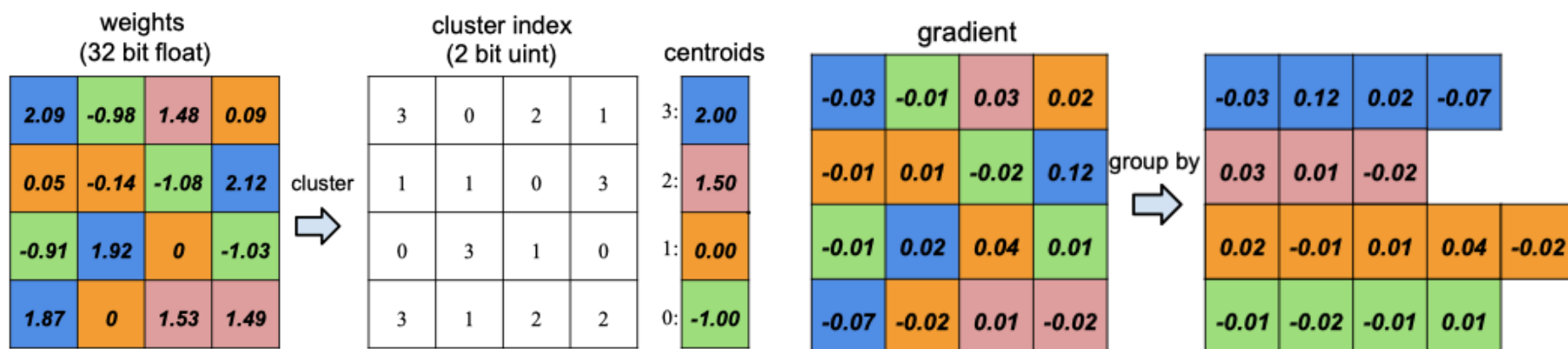
	Inference	Training
Algorithm	Pruning : P-Threshold, +Training Weight Sharing : <ul style="list-style-type: none"> w-KNN → Quant, Huffman coding Quantization : Floating P. → Fixed P. Low Rank Approximation Binary/Ternary Net Computational Transform <ul style="list-style-type: none"> Winograd, FFT, Strassen 	Parallelization : <ul style="list-style-type: none"> Data/Parameter/Hyp.Par Mixed Precision MAC : FP16 + FP32 Model Distillation : <ul style="list-style-type: none"> Teachers(GGN,VGG,Res)&Student DSD(Dense-Sparse-Dense) Training Mixed Precision Net : <ul style="list-style-type: none"> BinaryConnect, QNN
Hardware	Architecture : <ul style="list-style-type: none"> Temporal(SIMD,SIMT with Global Buffer) : General purpose GPUs Spatial(Local Memory in Dist. Processing Element) Data Flow Taxonomy : <ul style="list-style-type: none"> Weight Stationary, Output Stationary No local reuse (M-Bottleneck: Roofline) Row stationary Data type : Int8/16 for inference(TPU1), Mix FP16/32 for training(TPU2) Data Reuse : Conv, Feature map, Filter Kernel Computation (Toeplitz Matrix, TensorCore, Matrix-Multiply Unit) Sparse, Compressed Model (EIE, NVIDIA)	

Inference Acceleration Algorithms

Pruning

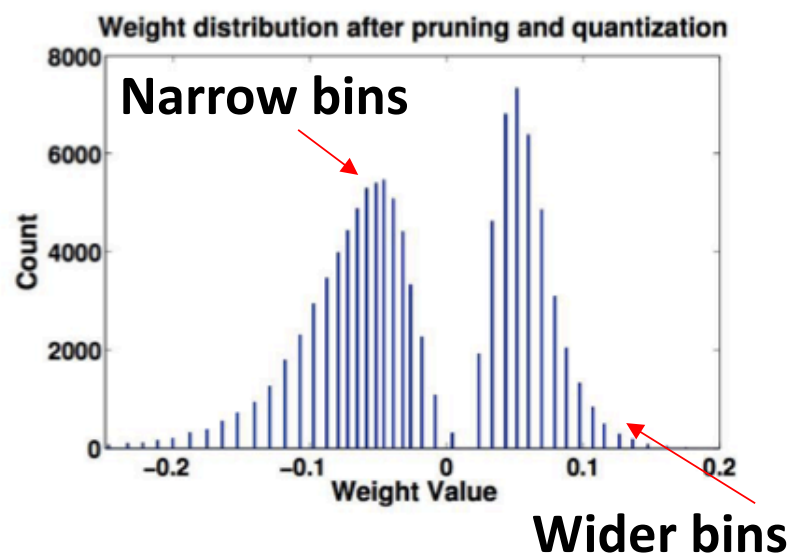
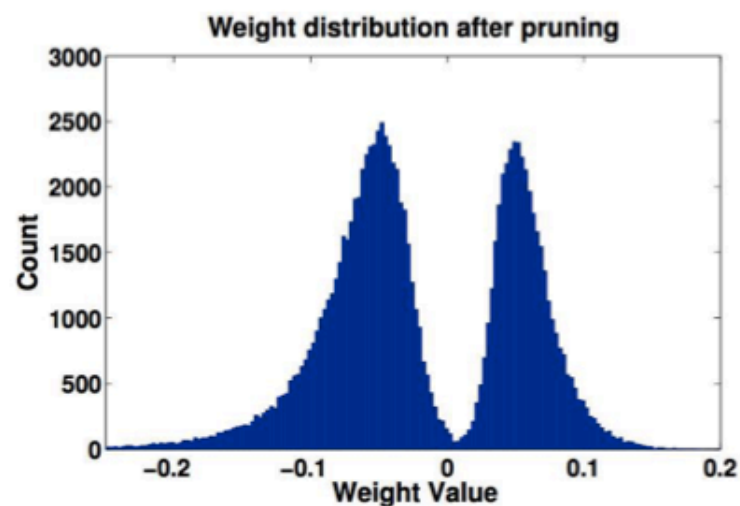


Weight Sharing

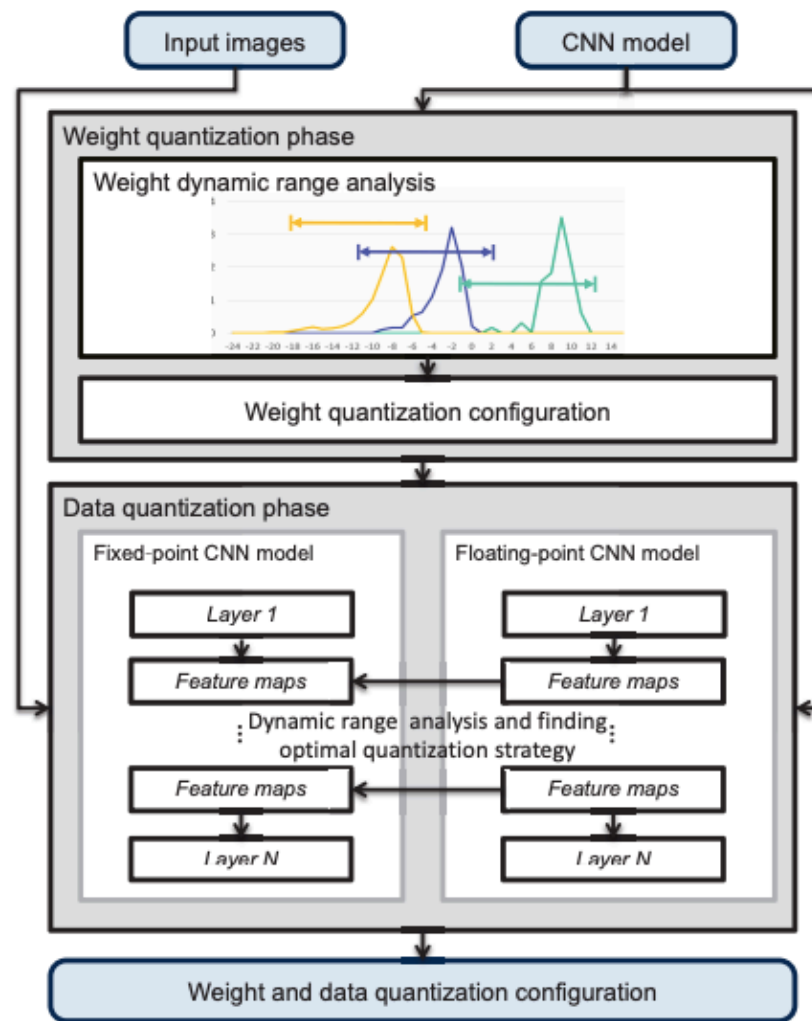


Inference Acceleration Algorithms

Weight Sharing & Quantization

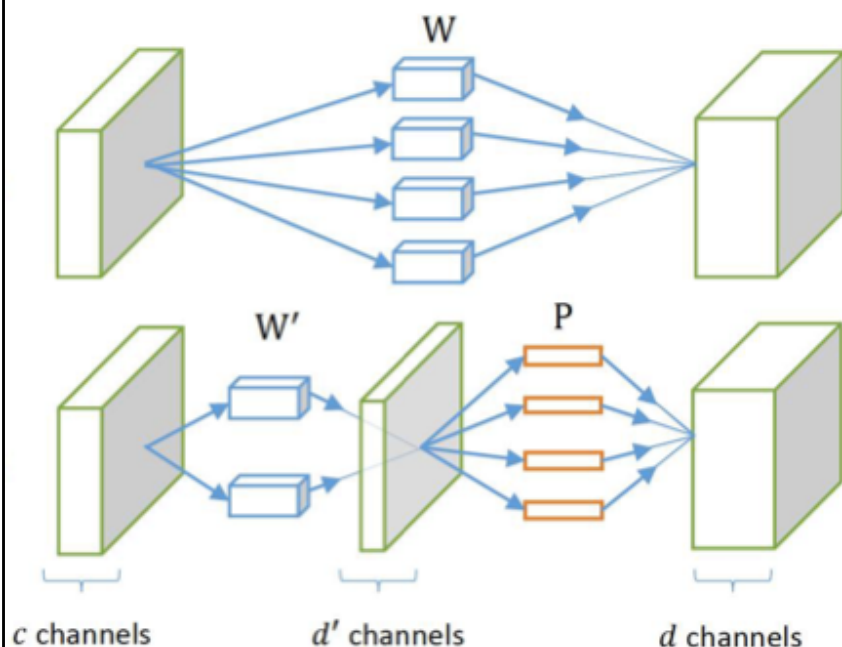


Quantization (Floating→Fixed)



Inference Acceleration Algorithms

Low rank approximation



Binary/Ternary Net for Inference(XOR-Net)

	Network Variations	
Standard Convolution	Real-Value Inputs 	Real-Value Weights
Binary Weight	Real-Value Inputs 	Binary Weights
BinaryWeight Binary Input (XNOR-Net)	Binary Inputs 	Binary Weights

Inference Acceleration Algorithms

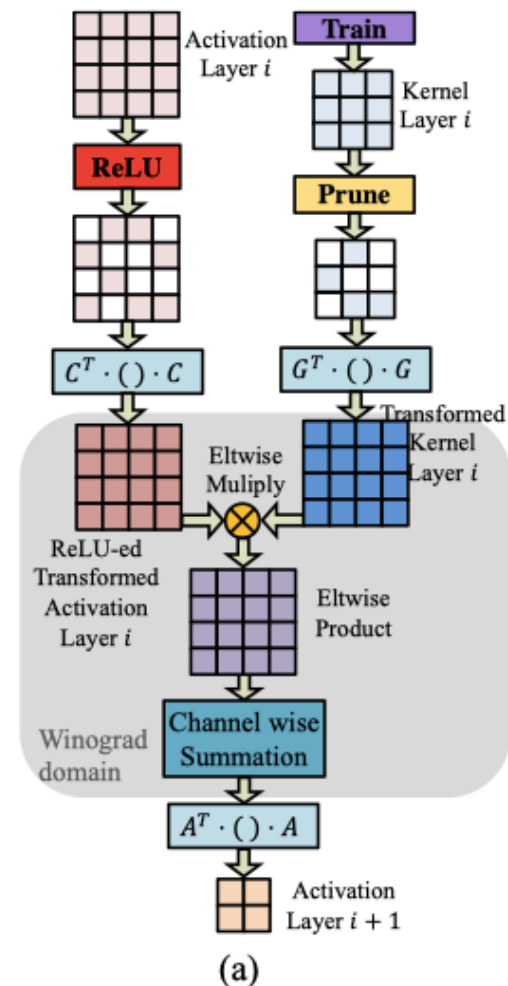
Winograd Transform

$$\begin{array}{|c|c|c|} \hline \text{Filter} & & \\ \hline g_{00} & g_{01} & g_{02} \\ \hline g_{10} & g_{11} & g_{12} \\ \hline g_{20} & g_{21} & g_{22} \\ \hline \end{array} * \begin{array}{|c|c|c|c|} \hline \text{Input Fmap} & & & \\ \hline d_{00} & d_{01} & d_{02} & d_{03} \\ \hline d_{10} & d_{11} & d_{12} & d_{13} \\ \hline d_{20} & d_{21} & d_{22} & d_{23} \\ \hline d_{30} & d_{31} & d_{32} & d_{33} \\ \hline \end{array} = \begin{array}{|c|c|} \hline \text{Output Fmap} & & \\ \hline y_{00} & y_{01} \\ \hline y_{10} & y_{11} \\ \hline \end{array}$$

3x3 filter & 4x4 inputs
= 9x4 = **36 multiplications**

Winograd Transform
→ 4x4 filter & 4x4 inputs
= **16 element-wise multiplications**

Winograd Transform & Prune



A. Lavin et. al., "Fast Algorithms for Convolutional Neural Networks", arXiv 2015

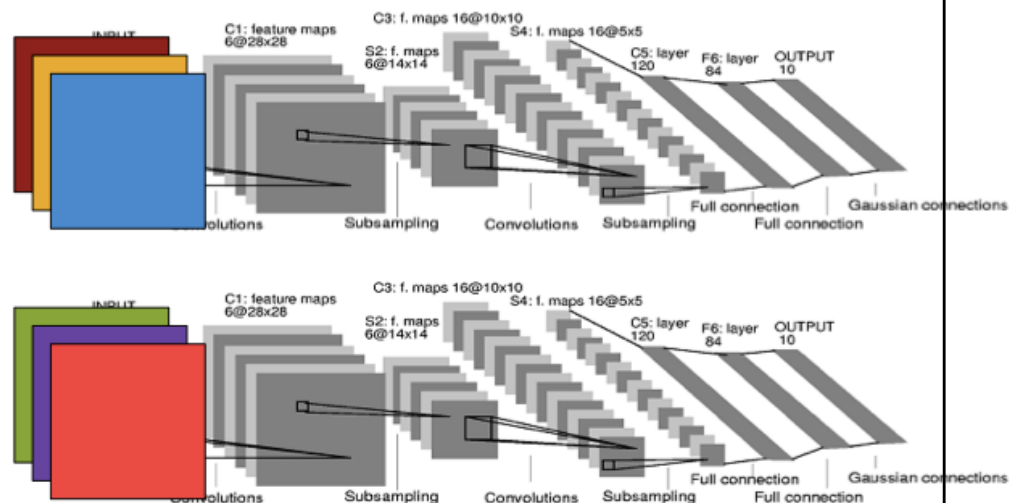
Liu et al. Figure "Efficient Sparse-Winograd Winograd Convolutional convolution with Neural sparse Networks", ICLR 2017

Strategies for DL Acceleration

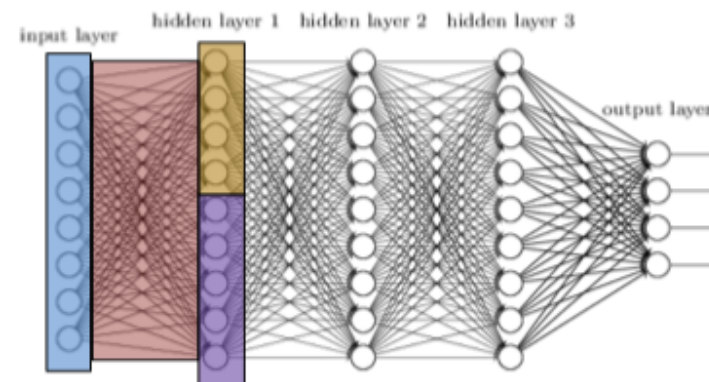
	Inference	Training
Algorithm	Pruning : P-Threshold, +Training Weight Sharing : <ul style="list-style-type: none"> w-KNN → Quant, Huffman coding Quantization : Floating P. → Fixed P. Low Rank Approximation Binary/Ternary Net Computational Transform <ul style="list-style-type: none"> Winograd, FFT, Strassen 	Parallelization : <ul style="list-style-type: none"> Data/Parameter/Hyp.Par Mixed Precision MAC : FP16 + FP32 Model Distillation : <ul style="list-style-type: none"> Teachers(GGN,VGG,Res)&Student DSD(Dense-Sparse-Dense) Training Mixed Precision Net : <ul style="list-style-type: none"> BinaryConnect, QNN
Hardware	Architecture : <ul style="list-style-type: none"> Temporal(SIMD,SIMT with Global Buffer) : General purpose GPUs Spatial(Local Memory in Dist. Processing Element) Data Flow Taxonomy : <ul style="list-style-type: none"> Weight Stationary, Output Stationary No local reuse (M-Bottleneck: Roofline) Row stationary Data type : Int8/16 for inference(TPU1), Mix FP16/32 for training(TPU2) Data Reuse : Conv, Feature map, Filter Kernel Computation (Toeplitz Matrix, TensorCore, Matrix-Multiply Unit) Sparse, Compressed Model (EIE, NVIDIA)	

Training Acceleration Algorithms

Data Parallelization



FC Model Parallelization

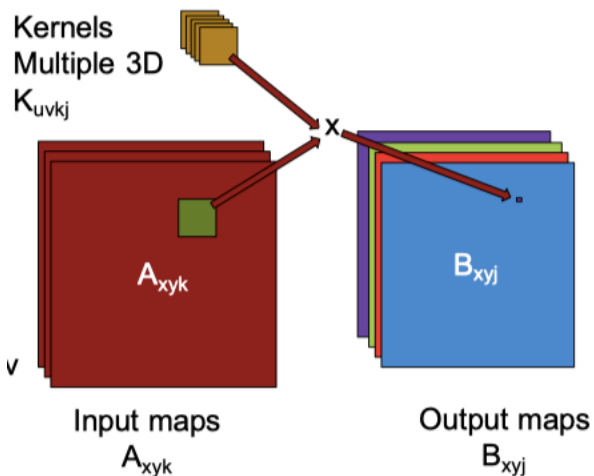


$$\begin{bmatrix} b_i \\ b_i \end{bmatrix} = \begin{bmatrix} W_{ij} \\ W_{ij} \end{bmatrix} \times \begin{bmatrix} a_j \end{bmatrix}$$

weight matrix

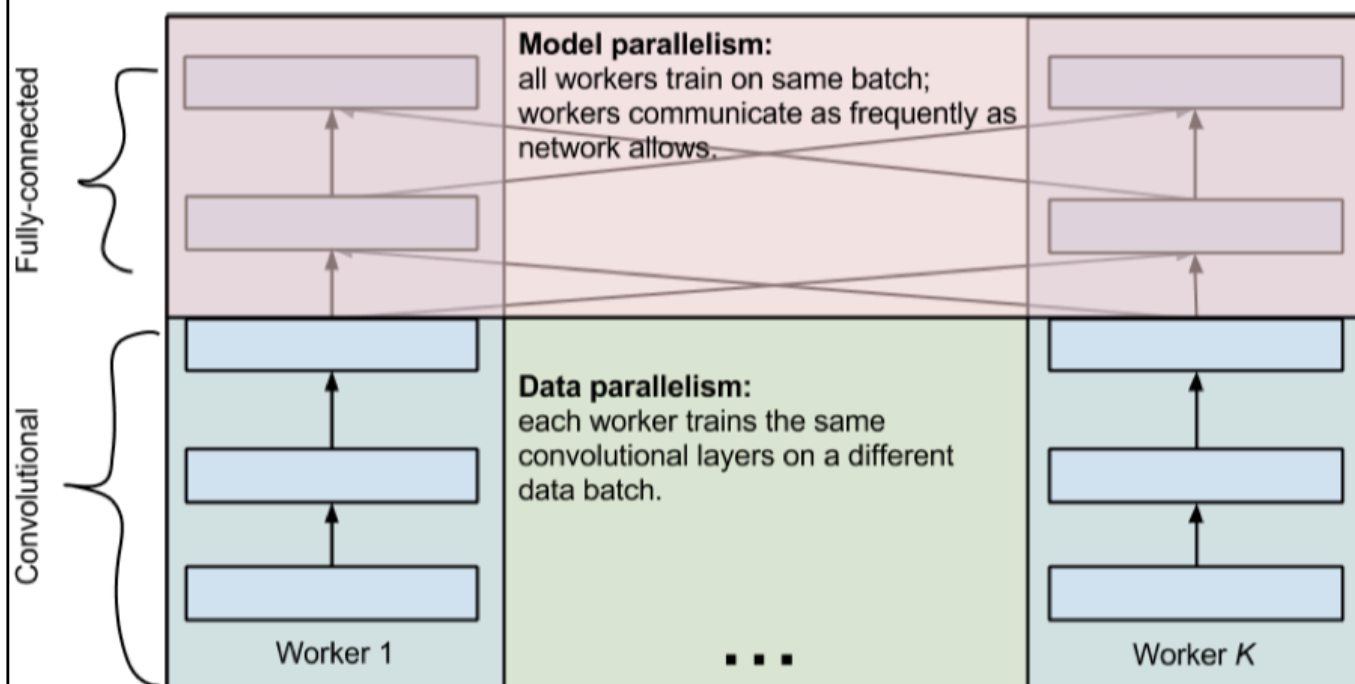
Output activations Input activations

Model Parallelization



Training Acceleration Algorithms

Alex's Weird CNN Parallelization Trick

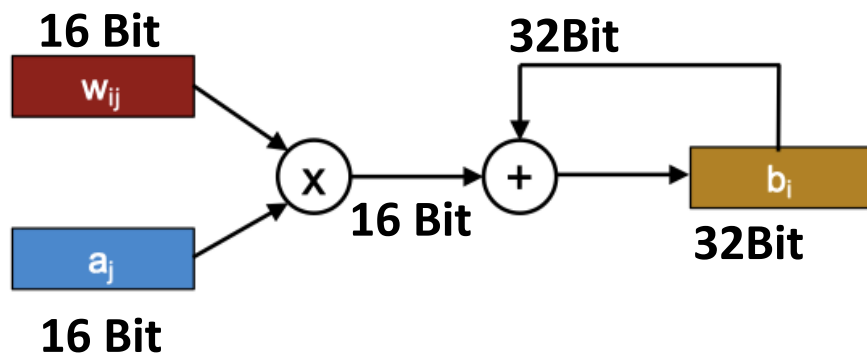


Metrics	LeNet-5	AlexNet
Top-5 error	n/a	16.4
Input Size	28x28	227x227
# of CONV Layers	2	5
Filter Sizes	5	3, 5, 11
# of Channels	1, 6	3 - 256
# of Filters	6, 16	96 - 384
Stride	1	1, 4
# of Weights	2.6k	2.3M
# of MACs	283k	666M
# of FC layers	2	3
# of Weights	58k	58.6M
# of MACs	58k	58.6M
Total Weights	60k	61M
Total MACs	341k	724M

- Most of Parameters → FC
- Most of MACs → Filter

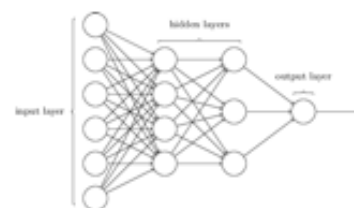
Training Acceleration Algorithms

Mixed Precision MAC



Mixed Precision MAC

Low Precision Net.

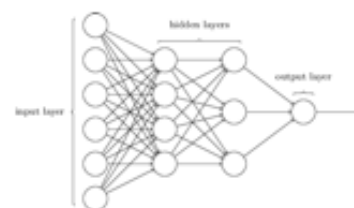


Forward/Backward Prop : $o(x)$ & $\frac{\partial \delta}{\partial w_b}$

Recall
analog weights (w)
for update

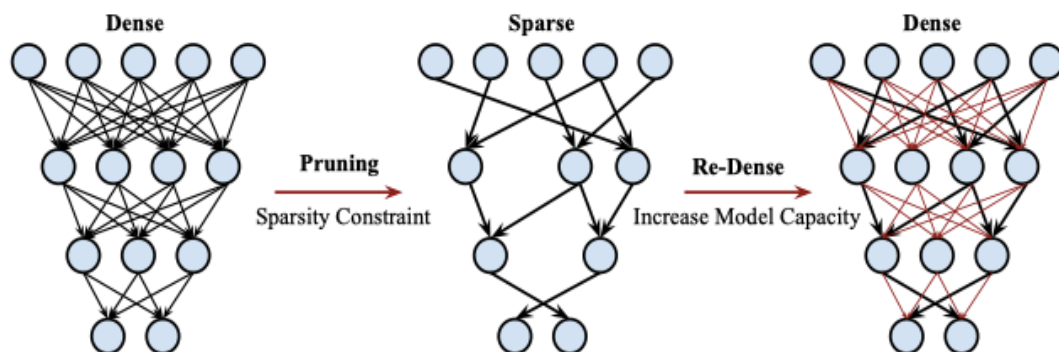
Stochastic
Binarization
 $w \rightarrow w_b$

High Precision Net.



Parameter update : $w' = w - \lambda \frac{\partial \delta}{\partial w}$

Dense-Sparse-Dense Training

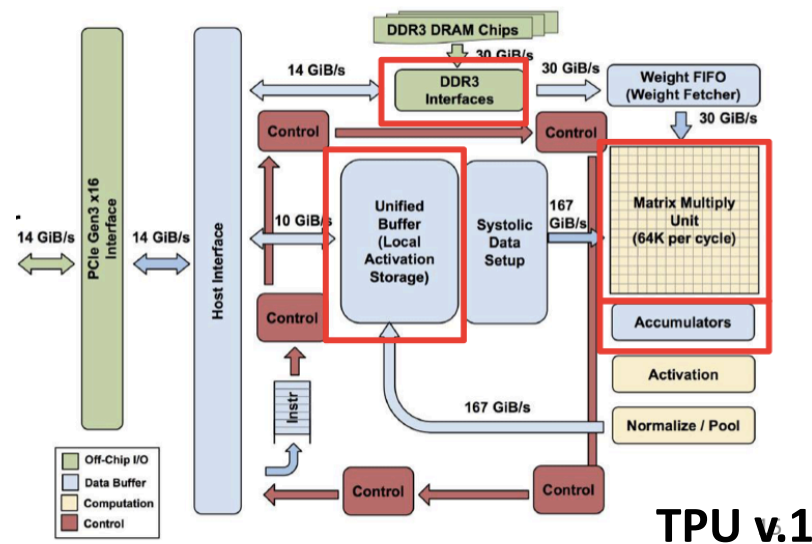
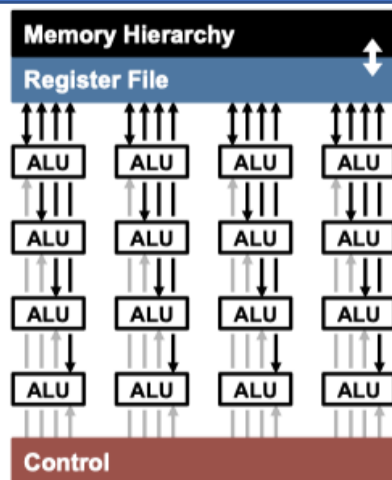


Strategies for DL Acceleration

	Inference	Training
Algorithm	<p>Pruning : P-Threshold, +Training</p> <p>Weight Sharing:</p> <ul style="list-style-type: none"> w-KNN → Quant, Huffman coding <p>Quantization: Floating P. → Fixed P.</p> <p>Low Rank Approximation</p> <p>Binary/Ternary Net</p> <p>Computational Transform</p> <ul style="list-style-type: none"> Winograd, FFT, Strassen 	<p>Parallelization:</p> <ul style="list-style-type: none"> Data/Parameter/Hyp.Par <p>Mixed Precision MAC: FP16 + FP32</p> <p>Model Distillation:</p> <ul style="list-style-type: none"> Teachers(GGN,VGG,Res)&Student <p>DSD(Dense-Sparse-Dense) Training</p> <p>Mixed Precision Net:</p> <ul style="list-style-type: none"> BinaryConnect, QNN
Hardware	<p>Architecture:</p> <ul style="list-style-type: none"> Temporal(SIMD,SIMT with Global Buffer) : General purpose GPUs Spatial(Local Memory in Dist. Processing Element) <p>Data Flow Taxonomy:</p> <ul style="list-style-type: none"> Weight Stationary, Output Stationary No local reuse (M-Bottleneck: Roofline) Row stationary <p>Data type: Int8/16 for inference(TPU1), Mix FP16/32 for training(TPU2)</p> <p>Data Reuse: Conv, Feature map, Filter</p> <p>Kernel Computation(Toeplitz Matrix, TensorCore, Matrix-Multiply Unit)</p> <p>Sparse, Compressed Model(EIE, NVIDIA)</p>	

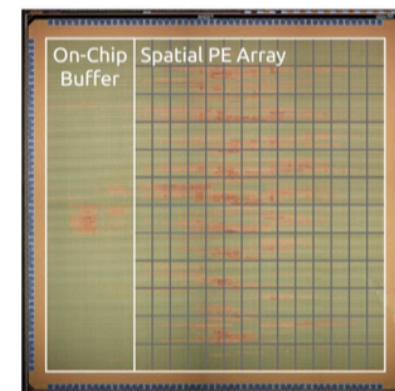
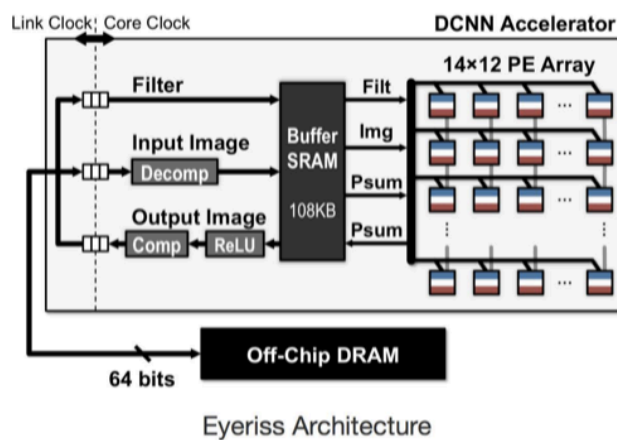
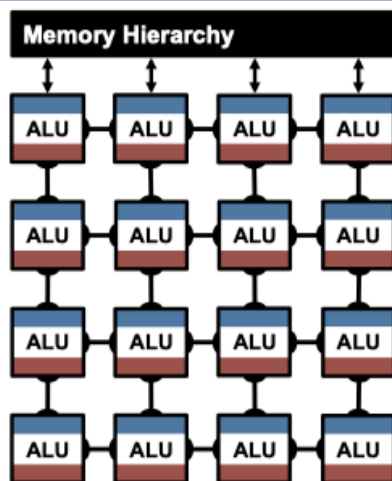
Hardware DL Acceleration

Temporal Architectures (SIMD, SIMT)



TPU v.1

Spatial Architectures (Local Memory & PE)



Die Photo

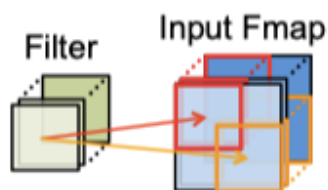
MIT:Eyeriss

Hardware DL Acceleration

Data Reuse

Convolutional Reuse

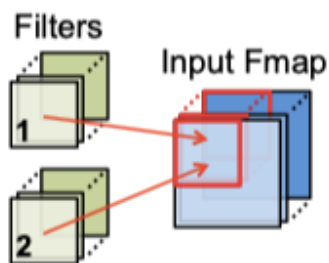
CONV layers only
(sliding window)



Reuse: **Activations**
Filter weights

Fmap Reuse

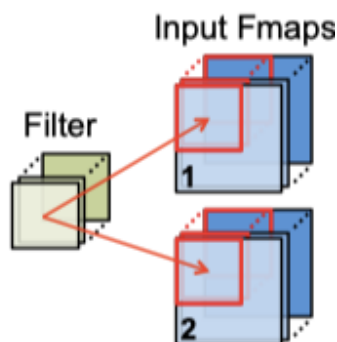
CONV and FC layers



Reuse: **Activations**

Filter Reuse

CONV and FC layers
(batch size > 1)



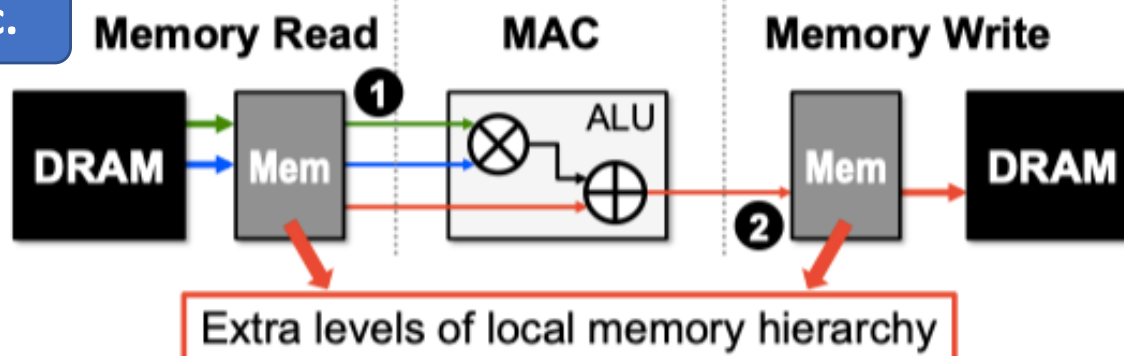
Reuse: **Filter weights**

Exploiting Sparsity

$$\vec{a} \begin{pmatrix} 0 & a_1 & 0 & a_3 \end{pmatrix} \times \begin{matrix} PE0 \\ PE1 \\ PE2 \\ PE3 \end{matrix} \begin{pmatrix} w_{0,0} & w_{0,1} & 0 & w_{0,3} \\ 0 & 0 & w_{1,2} & 0 \\ 0 & w_{2,1} & 0 & w_{2,3} \\ 0 & 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} b_0 \\ b_1 \\ -b_2 \\ b_3 \\ -b_4 \\ b_5 \\ b_6 \\ -b_7 \end{pmatrix} \xRightarrow{ReLU} \vec{b} \begin{pmatrix} b_0 \\ b_1 \\ 0 \\ b_3 \\ 0 \\ b_5 \\ b_6 \\ 0 \end{pmatrix}$$

Virtual Weight	$w_{0,0}$	$w_{0,1}$	$w_{4,2}$	$w_{0,3}$	$w_{4,3}$
Relative Index	0	1	2	0	0
Column Pointer	0	1	2	3	

Reuse & Local Acc.

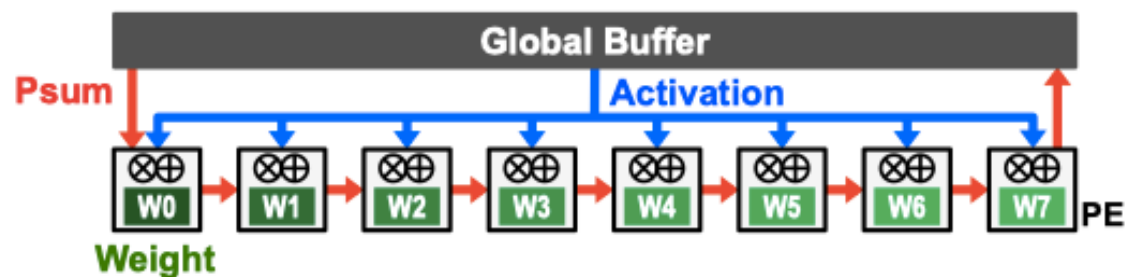


Opportunities: **1** data reuse **2** local accumulation

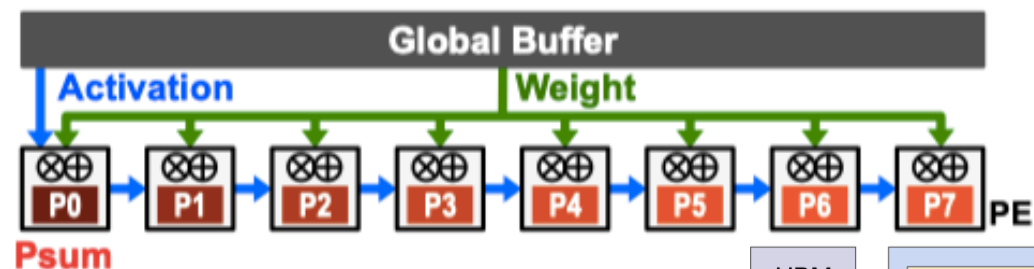
Hardware DL Acceleration

Data Flow Taxonomy

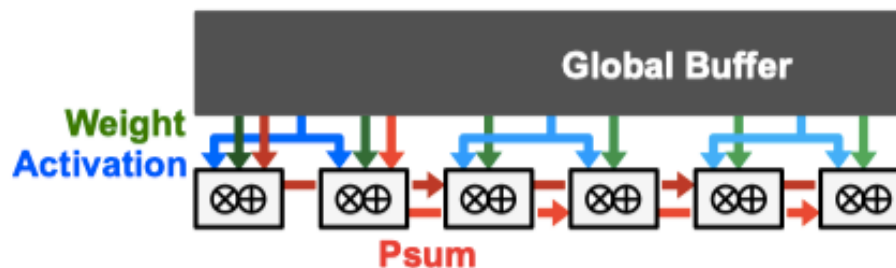
Weight
Stationary



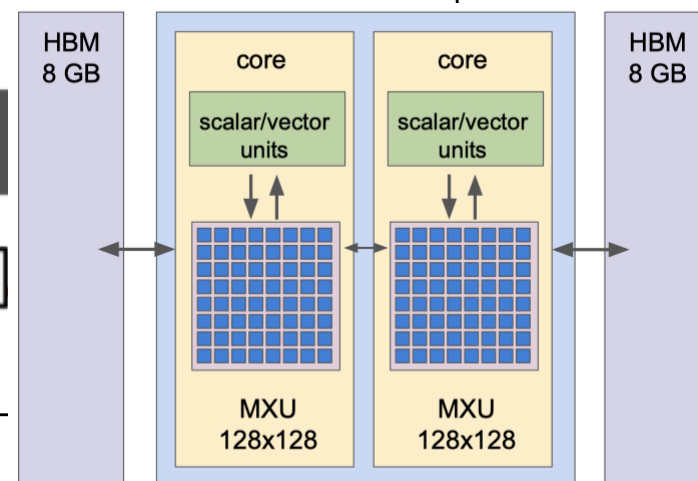
Output
Stationary



No Local
Reuse
(TPU1)



(TPU2)



Hardware DL Acceleration

Kernal Computation

Filter * Input Fmap = Output Fmap

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} * \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

Toeplitz Matrix
(w/ redundant data)

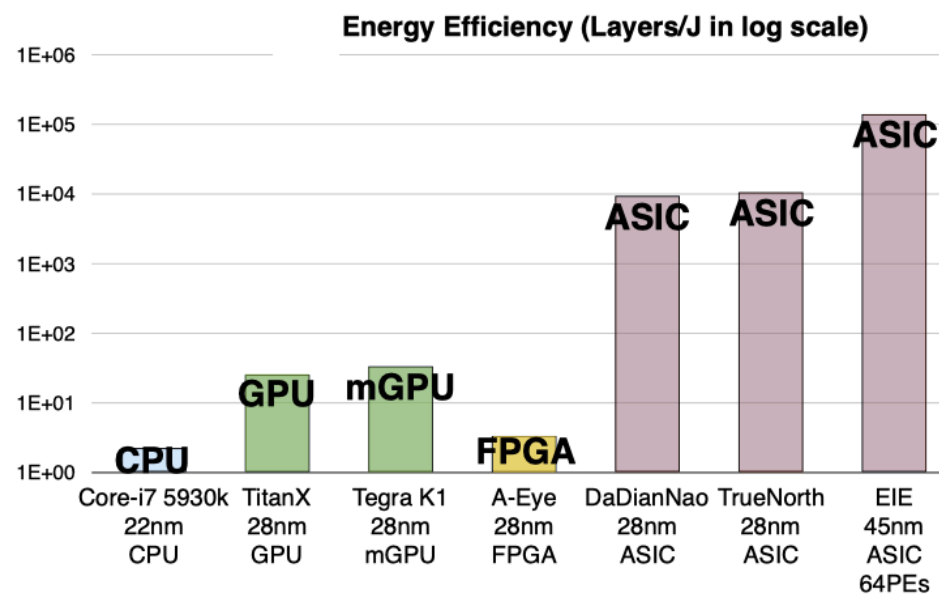
$$\begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 4 & 5 \\ 2 & 3 & 5 & 6 \\ 4 & 5 & 7 & 8 \\ 5 & 6 & 8 & 9 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$$

Nvidia's Tensor Core

$$\begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 FP16 FP16 or FP32

Energy Efficiency



Strategies for DL Acceleration

	Inference	Training
Algorithm	<p>Pruning : P-Threshold, +Training</p> <p>Weight Sharing:</p> <ul style="list-style-type: none"> w-KNN → Quant, Huffman coding <p>Quantization: Floating P. → Fixed P.</p> <p>Low Rank Approximation</p> <p>Binary/Ternary Net</p> <p>Computational Transform</p> <ul style="list-style-type: none"> Winograd, FFT, Strassen 	<p>Parallelization:</p> <ul style="list-style-type: none"> Data/Parameter/Hyp.Par <p>Mixed Precision MAC: FP16 + FP32</p> <p>Model Distillation:</p> <ul style="list-style-type: none"> Teachers(GGN,VGG,Res)&Student <p>DSD(Dense-Sparse-Dense) Training</p> <p>Mixed Precision Net:</p> <ul style="list-style-type: none"> BinaryConnect, QNN
Hardware	<p>Architecture:</p> <ul style="list-style-type: none"> Temporal(SIMD,SIMT with Global Buffer) : General purpose GPUs Spatial(Local Memory in Dist. Processing Element) <p>Data Flow Taxonomy:</p> <ul style="list-style-type: none"> Weight Stationary, Output Stationary No local reuse (M-Bottleneck: Roofline) Row stationary <p>Data type: Int8/16 for inference(TPU1), Mix FP16/32 for training(TPU2)</p> <p>Data Reuse: Conv, Feature map, Filter</p> <p>Kernel Computation(Toeplitz Matrix, TensorCore, Matrix-Multiply Unit)</p> <p>Sparse, Compressed Model(EIE, NVIDIA)</p>	